

# COMP90042 Project 2021: Rumour Detection and Analysis on Twitter

1164051

## Abstract

Rumour detection on Twitter is a significant issue. In the first part, the problem of rumour detection on twitter dataset addressed by applying a number of machine learning and deep learning models, such as SVM, C-LSTM and BERT, to perform rumour analysis on source tweets as well as reply tweets. In the second part, the model of trained BERT rumour classifier from the first task applied to a set of provided COVID-19 tweets dataset to detect rumours and analyse how topics, users and emotions of rumours differ from non-rumours.

## 1 Tweet Rumour Detection Introduction

Rumours spread rapidly via Twitter and can have a huge economic and social impact. In this report, I will attempt to conduct rumour detection on “tweets” using various machine learning algorithms. Commonly, there are two possible ways for the rumour detection: Linguistic Approach and Network Approach. For the former, the rumour content should be extracted based on language patterns. The misinformation would have some similar language patterns such as fear mongering. While for Network Approach, the metadata for the information can be taken into consideration, such as the user characters in this case. A tweet from an unreliable source (rumour spreader user) is likely to be a rumour. In this report, I will consider only the Linguistic Approach as described and identify rumours based on only text information (tweets “text”).

**Note:** All the figures and tables were shown in the Appendix section.

## 2 Methodology

### 2.1 Dataset Pre-processing

Twitter data must be normalized to create a dataset that can be easily learned by classifiers since tweets have certain special characteristics such as emotions, user mentions, URLs, etc. There are

several pre-processing steps performed to clean the raw data and reduce its size:

- Use the NLTK Tweet Tokenizer package to tokenize the tweets.
- Remove all the URLs by matching expression (`www\.[\S]+`)(`https?:\/\/[\S]+`).
- Remove all the user mentions by matching expression (`@[\s]+`).
- Convert the tweet to lower case.
- Remove single character words, non-English words and custom stopwords expanded from NLTK stopword list.
- Lemmatise all words by applying WordNetLemmatizer package tweet.

### 2.2 Feature Engineering

Most rumour spreaders would use certain language strategies to increase its credibility to the audience. Therefore, how the text features can capture the language patterns for rumours will directly impact the accuracy of the whole rumour detection system. Here in my implementation, I applied three type of word features representation: BOW, TF-IDF and word sequence in different models.

**BOW:** BOW regards the post-processed words in a tweet text as separated with equally important tokens. This representation method can give much information on the text itself, as in our problem the words in rumours much likely to express an unobjective emotion or negative case.

**TF-IDF:** I also used TF-IDF in BOW in SVM models, however, this method can only represent simple word mentions in texts and unable to capture syntax or styles of rumours. Rumours and non-rumours may share the similarity of tokens, but the different language syntax and structure will lead to totally different meaning.

**Word Sequence:** Sequence respecting approaches have an edge over BOW when a synthetic tweet corpus dominated by sentences and consideration to the position of the words are important. When I trained some neural network models, I applied some deep learning approaches such as LSTM,

GNU to model the tweets as embedding string of words.

### 2.3 Models

Three models are used in the project to create the detection system, one for baseline model using binary class classification SVM, two for model improvements that one is neural network mix model CLSTM, another is BERT. "train.json" dataset will be used in training process and hyper parameter tuning and evaluation will be conducted on "dev.json" dataset.

Note: I divided tweets in both train and dev datasets into source tweets and source + reply tweets, training and evaluating respectively.

**SVM:** SVM is a simple non-probabilistic binary linear classifier. To compare the data features performance between BOW and TF-IDF, I built the model with both two features from post-processed tweet data and implemented by using SVC classification from scikit-learn SVM package in Python with fine-tuned radial basis function kernel where  $C = 1.0$ .

**C-LSTM:** I used keras with TensorFlow backend to implement the C-LSTM Neural Network model, which combined the strengths of RNN and CNN architectures for sentence representation and text classification. C-LSTM utilizes CNN to extract a sequence of higher-level phrase from source tweets, then the vectors generated from the CNN are fed into a LSTM recurrent neural network for rumour detection for tweets. The details illustrate in Figure 2.3.1. In my implementation, after initializing the first layer embeddings of sequential inputs, the output from the CNN layer (second layer) was directly concatenated with LSTM layers (third layer).

**Fine-Tuning BERT:** Thinking about the tweets data between the source and reply will have a contextual relationship, as well as a time-based relationship after sorted by tweets created time, I used the BERT BASE since it is a bi-direction trained transformer language model which contains an attention mechanism to learn the contextual relationship between words and then generate deep bidirectional context representations by jointly conditioning on both left and right context in all layers. In my implementation (Figure 2.3.2), after feeding the input sequence to BERT

and pass the contextualised embedding of [CLS] produced by BERT to a simple feedforward network classifier, I add one additional output layer with a SoftMax function for rumour classification, so that it can produce a single scalar value to denote the probability of the rumour class. I also applied fine-tuning in BERT that extends the training max length to 256 since the average input sentences distributed between 150 to 200.

### 3 Evaluation and Analysis

**SVM:** Based on the results (Table 3.1), SVM performs badly and fails due to very low recall across all classes and results in low F1 scores (Average 69%), however, it performs well in precision above 0.8, indicating that the model classifies a few negative cases as rumour, it shows that SVC classifier tends to predict a label as positive. We can see TF-IDF representation recorded the best performance in source tweets, the reason I supposed that TF-IDF will filter out irrelevant (noise) words and with training only on the source tweets, it has less dimensions to expend which results in higher efficient since SVM has a poor explanatory power for high-dimensional mappings of kernel functions, especially radial basis functions.

**C-LASTM:** I have applied CNN and RNN, C-LSTM models separately to compare the performances. According to the results (Table 3.1), they all have a significant improvement on F1 score (Average 78%) compared with SVM, though it still not up to expectation. C-LSTM performs best on F1 score, where CNN had a slightly lower score following the worst performance on RNN. This can be explained by that sequential processing over time works well on whole tweets, with combined the advantages of features extraction from CNN and long context memory from LSTM. However, the model performance is unstable, and computation and time consumptions are not efficient.

**BERT:** BERT performs consistently excellent for each class (Table 3.1) and reach to the highest F1 score at 81.73% on the whole dataset (source + reply tweets). That means, BERT can adjust the weights associated with the model to better represent text. In more specific, during classifier fine-tuning, the starting points of the weights are closer to values that correctly model Twitter data.

However, there is a huge time consumption during the training process because of the slow converges on model since only 15% of the data in each batch size is involved in the prediction.

**Comparison Between Models:** Comparing the best performance between the three models (Table 3.2), BERT has an obviously higher F1 score and more stable performance than the others. It proves again that a Bi-direction trained language model has a deeper understanding of the context than a one-way language model, which becomes to the most suitable model to process the tweets dataset.

## 4 CodaLab Competition

I finally used an ensemble method taking a majority vote over the predictions of the best performance in the above 3 different models achieving an accuracy of 83.39%. (Table 4)

## 5 Covid-19 tweets Analysis

### 5.1 Introduction

Since late February 2020, the COVID-19 pandemic has come to dominate both traditional news and social media platforms, and misinformation such as fake news, conspiracy theories and rumours thrive during these uncertain times. The aim of this section is to figure out what kinds of COVID-19 rumours are being distributed on Twitter. After applying the best performing rumour detection model from part 1 to the given Covid-19 based tweets dataset, projected labelled rumour data could be obtained to distinguish rumours from non-rumours in this COVID-19 data and used to analyse the users' characteristics and tweets texts. In this dataset, the findings show a variety of fascinating observations about users, topics, and emotions. For example, rumour-spreaders, usually have low follower but a high following count and they prefer to discuss politics (mostly party blaming), are more emotionally charged (e.g., anger) with more negative sentiments.

### 5.2 Data design

After applying the BERT rumour detection on the total 17458 events in COVID-19 dataset, there are 1548 events were labelled as "rumour" and 15910 events were labelled as "non-rumour". By counting the retweet and tweet counts on each day, we can visualize the volume of non-rumour and

rumour tweets over time respectively as the Appendix section 2, Figures 5.2.1 and 5.2.2 shown. We can see that they both have some traffic of COVID-19 related tweets from late January 2020, although they do not really pick up until early-March and I suspect the spike of activity may be due the World Health Organisation declaring it as a pandemic on 12th March 2020. There is a significant high point on the rumour figure at around the date of 24 April 2020, I supposed this may be due the rumour words from Trump that the COVID-19 can be treated by injecting disinfectant which has drawn much attention to discuss.

In terms of pre-processing, I use the NLTK TweetTokenizer package to tokenize the tweets, and lowercase and lemmatise all words by applying WordNetLemmatizer package, as well as eliminate digits, hyperlinks and @usernames. I also use an expanded NLTK stopword list to filter stopwords, which includes COVID-19-related stopwords like covid19 and coronavirus. For topic analysis, hashtags are also omitted.

### 5.3 Results and Analysis

**Topic analysis:** To find the popular topics discussed in datasets, I retrieved the top frequent unigram and bi-gram words which extracted from source tweets (here I did not add reply tweets since in some case, reply tweet might interference with the central topics of a tweet event) and then visualized the frequency dictionaries in Wordcloud format as well as line charts (Appendix section 2 Figures 5.3.1~5.3.8). After comparison, bigram wordclouds (Figures 5.3.1, 5.3.3) show more reasonable results and we can see several broad topics: (1) COVID19-status reports (tested positive, confirmed case, new case, report new, death toll); (2) health advice (social distancing, public health, and wear mask) appeared frequently in both two classified datasets where the topic of US politics (president trump and white house) accounts for the largest proportion for non-rumours, and COVID19-status reports for rumours. To better understand the topical difference between the two datasets over time changes, I tracked both bigram dictionaries sorted by months from January to July and Table 5.3.1 gives an overview of the top frequent topics over a month. As the pandemic grew, there is no surprise that the emerged keywords of rumours associated starting with "Wuhan epidemic" to "new cases", "death tolls" and "pandemic began" in the end.

Looking at non-rumours, the topics are very different: they are mostly related to “president trump” and health advice and turn to focusing on covid cases status in the late period. The reason I supposed that Trump’s highly active Twitter rumour remarks aroused heated discussion and dissatisfaction among people.

**Hashtag analysis:** To find the most popular hashtags of COVID-19 rumours and non-rumours, I separately extracted the hashtag words (start from the first character “#”) from the rumour and non-rumour raw source and reply tweets and filtered COVID-19-related hashtags like “#covid19”. After sorting the filtered hashtags frequencies and representing the top 20 popular hashtags of non-rumours and rumours as bar charts (Figure 5.3.8, 5.3.10). An Interesting finding is that the non-rumours figure shows a similarity with the above popular topics as the president trump still be a highly topical theme associated with many sarcasm and accusation emotions hashtags like “#TrumpVirus”, “#TrumpLiesPeopleDie”. For rumours, most popular hashtags related to the healthcare such as personal protection equipment and, unsurprisingly, “#WuhanVirus” and “#China” are also labelled as rumour hashtags.

**User analysis:** I focus only on users who published the source tweets in this analysis. Table 5.3.2 presents some statistics of these users for rumours and non-rumours and figure 5.3.11 visualizes the comparison of some items of rumour and non-rumour users’ characters. Interestingly, users who are involved in rumour creation tend to tweet more (higher post counts) but have fewer posts that are tagged as a favourite compared with non-rumour users. Besides, the rumour users have slightly fewer followers as well as the followings. Their account is also generally older than non-rumour users that the average created date is on 2010-12-01, whereas the non-rumours user accounts created at about 2011-04-30. I also retrieve the most mentioned users in raw tweets and selected top 20 @Username in rumours and non-rumours, respectively (Figures 5.3.12, 5.3.13). Except the @realDonaldTrump been to the most popular mentioned user in both classifications, we can see that some of well-known rumour news spreader like @FoxBusiness and @CNN were tagged as rumour accounts, on the other hand, @JoeBiden become an obvious non-rumour user which I suspect it is because compared with Trump's

unreliable policies, Biden's governance received more popular support during the pandemic.

**Sentiment/Emotion Analysis:** I choose an emotion prediction system to classify the emotion of tweets in the data to better recognize the public's sentiments during the COVID19 crisis. DeepMoji [1], a Bi-LSTM with an attention model trained on a large number of emoji occurrences in tweets, is used throughout this experiment and I labelled the source tweets data with 63 predefined emojis by applying their pre-trained model. Figure 5.3.14 illustrates the distribution of emojis for source tweets in rumours and non-rumours. Compared the two-emotion pie charts, we see a similar distribution for the top 5 emotions (👍, 😞, 😡, 🙏, 🤔), but a bit confused observation here is the top emoji for rumours is 👍 account for 9.07% whereas 😞 as the top emoji for non-rumours accounts for 8.88%, though the percent difference is less severe. Overall, based on the emoji represents, the emotions of the public generally express the anger and disagreement, the rest of them usually express encouragement, as well as a pray from the source tweets observation.

To accurate the sentiments of the public from tweets, I applied another sentiment detected model, package Baidu Senta Corpus [2], a type of sentiment Analysis system to interpretate and classify the sentiments (positive and negative) within text data. Figure 5.3.15 reveals the public sentiments towards source tweets and we can see that negative sentiment dominates both rumours and non-rumours, but substantially more in rumours than non-rumours (68% vs. 62%).

## 6 Conclusion

In this experiment, after implementing several ML algorithms: SVM, C-LSTM, BERT to build the rumour detection system, I found BERT outperforms the other models on tweet dataset which proved that bi-direction transformer more fit in processing tweets data. I finally reached the F1 score of 83.98% on CodaLab leader board via majority voting from 3 best predictions of the above 3 models. I applied BERT model to classify COVID-19-related rumours on given tweet datasets and provided a quantitative test to demonstrate analysis of rumours vs. non-rumours users, topics, and emotions, and discovered a number of insights.

## References

- [1] B. Felbo, A. Mislove, A. Søgaard, I. Rahwan, S. Lehmann, Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Copenhagen, Denmark, 2017, pp. 1615–1625.
- [2] Tian, H., Gao, C., Xiao, X., Liu, H., He, B., Wu, H., ... & Wu, F. (2020). SKEP: Sentiment knowledge enhanced pre-training for sentiment analysis. *arXiv preprint arXiv:2005.05635*.
- [3] Tian, L., Zhang, X., & Lau, J. H. (2020). # Democrats are destroying America: Rumour analysis on twitter during COVID-19. In CEUR Workshop Proc..Ashok K. Chandra, Dexter C. Kozen, and Larry J. Stockmeyer. 1981. [Alternation](#). *Journal of the* <https://doi.org/10.1145/322234.32224>.
- [4] Tian, L., Zhang, X., Wang, Y., & Liu, H. (2020, April). Early detection of rumours on twitter via stance transfer learning. In European Conference on Information Retrieval (pp. 575-588). Springer, Cham.
- [5] Zhou, C., Sun, C., Liu, Z., & Lau, F. (2015). A C-LSTM neural network for text classification. *arXiv preprint arXiv:1511.08630*.
- [4] Tian, L., Zhang, X., Wang, Y., & Liu, H. (2020, April). Early detection of rumours on twitter via stance transfer learning. In European Conference on Information Retrieval (pp. 575-588). Springer, Cham.

## Appendix

The figures referenced from the above report are shown below.

### Section 2

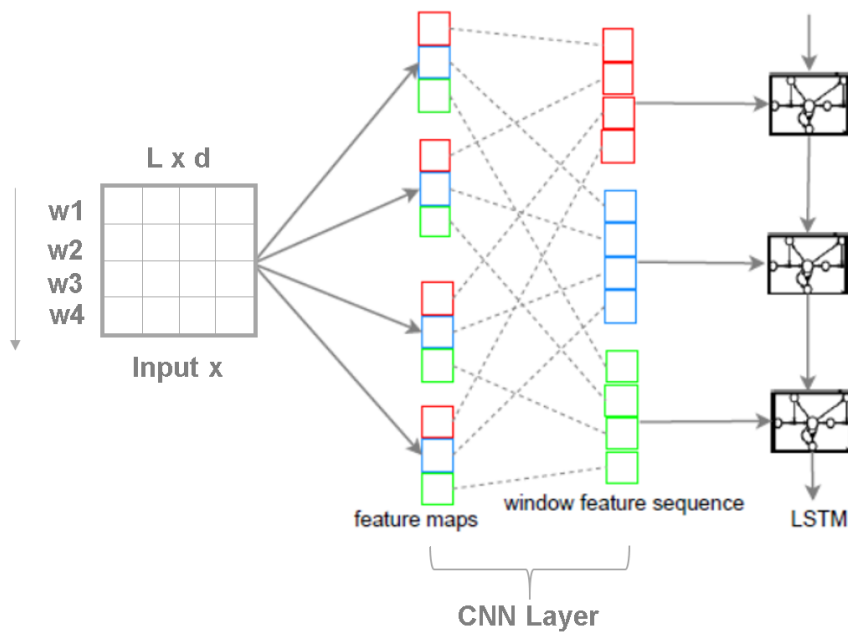


Figure 2.3.1. C-LSTM model description [5]

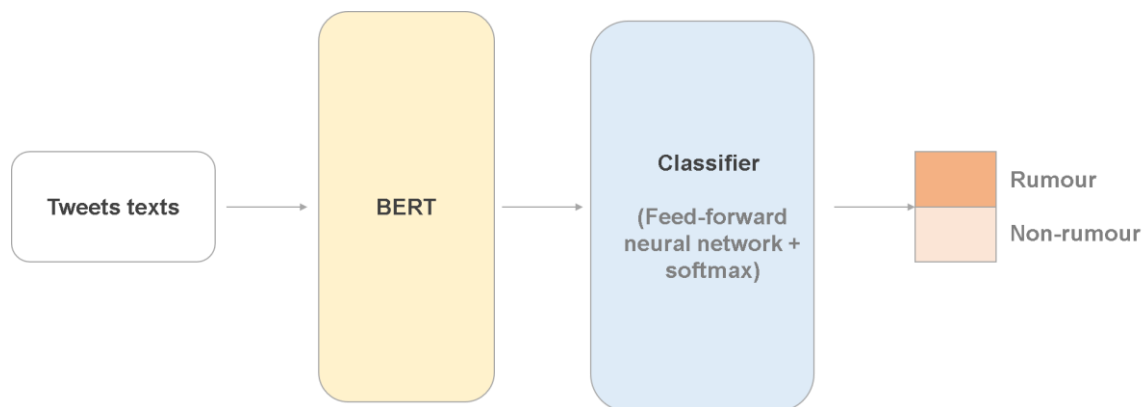
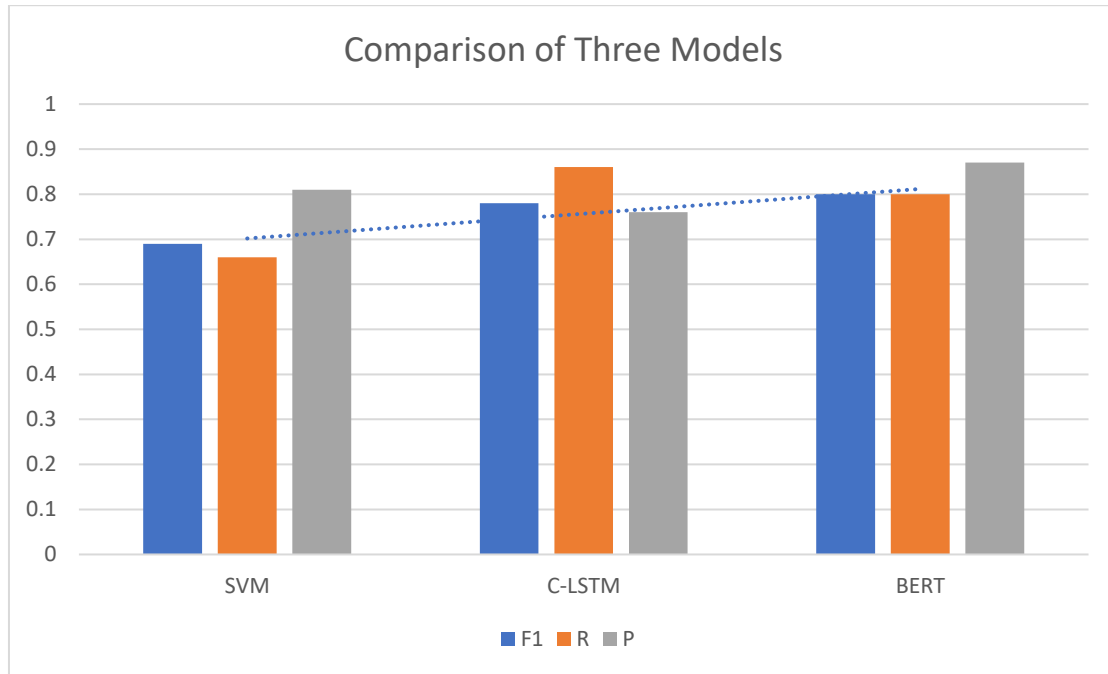


Figure 2.3.2. BERT model description

### Section 3

Model	Inputs	F1	R	P
SVM	TF-IDF (Source Tweets)	0.69	0.66	0.81
	TF-IDF (Whole Tweets)	0.67	0.66	0.79
	BOW (Source Tweets)	0.67	0.69	0.8
	BOW (Whole Tweets)	0.64	0.64	0.76
RNN	Source Tweets	0.75	0.67	0.84
CNN	Source Tweets	0.77	0.80	0.81
C-LSTM	Source Tweets	0.75	0.82	0.74
	Whole Tweets	0.78	0.86	0.76
BERT	Source Tweets	0.77	0.81	0.84
	Whole Tweets	0.80	0.80	0.87

**Table 3.1.** Results of Models



**Table 3.2.** Comparison of Three Models

#### Section 4

CodaLab Model	F1	R	P
Majority Vote			

**Table 4.** Results of CodaLab

#### Section 5.2

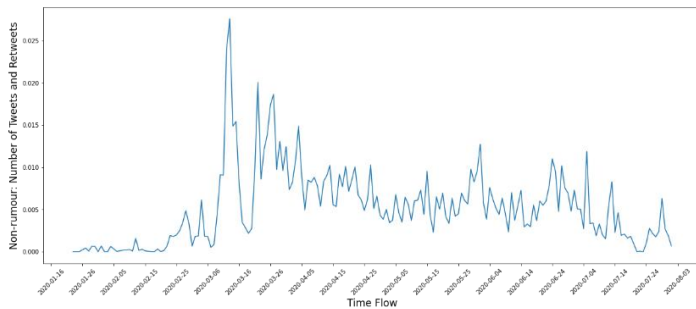


Figure 5.2.1. non-rumour tweets volume overtime.

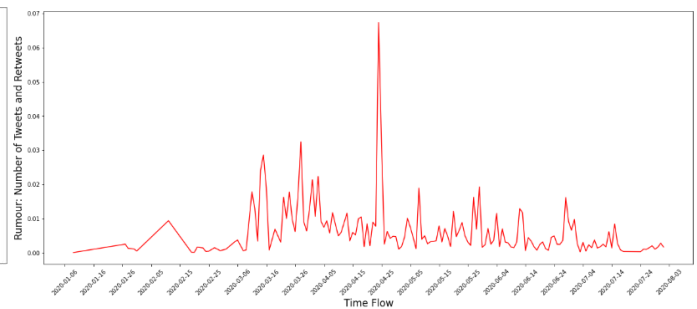


Figure 5.2.2: rumour tweets volume overtime.



The image displays two word clouds side-by-side, representing news headlines about the COVID-19 pandemic. The left word cloud, dated 2020-04-15, features prominent terms such as 'white house', 'social distancing', 'president trump', 'new case', 'tested positive', and 'task force'. The right word cloud, dated 2020-05-01, features prominent terms such as 'new case', 'death toll', 'confirmed case', 'reported new', and 'positive'. Both word clouds use a color palette of reds, oranges, and yellows, with word size indicating frequency or importance.

Figure 5.3.3. rumour bigram topics.

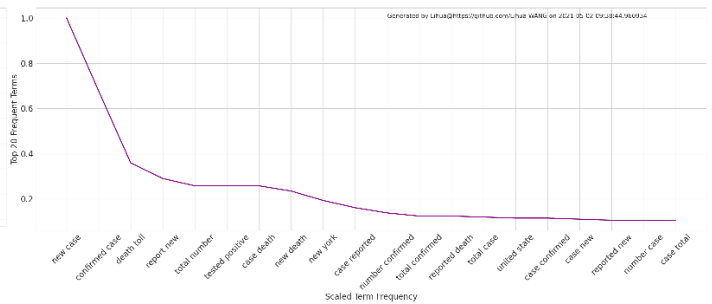


Figure 5.3.4. rumour bigram topics in line charts.

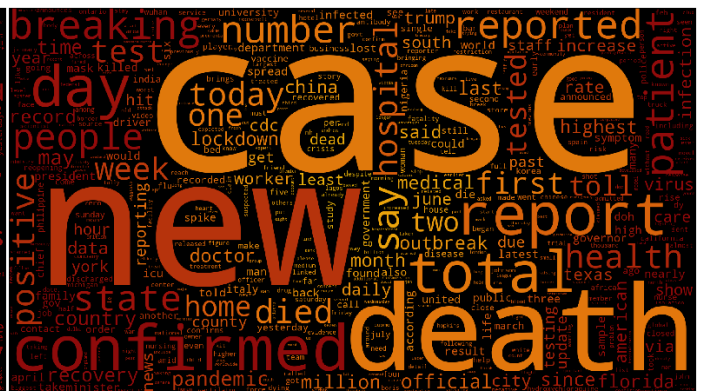


Figure 5.3.7. rumour unigram topics.

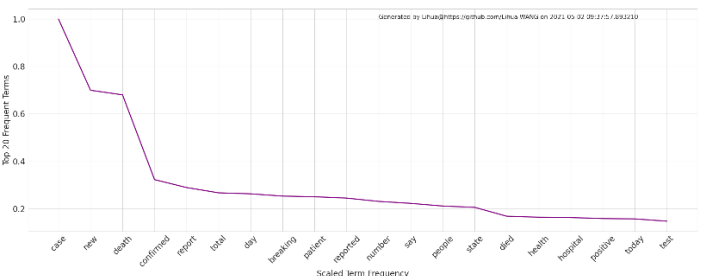


Figure 5.3.8. rumour unigram topics in line chart.

8



04	president trump, white house, tested positive, social distancing, stay home	death toll, new case, confirmed case, total number, case reported, case death, new york
05	president trump, tested positive, nursing home, new case, trump say, wear mask	new case, confirmed case, death toll, tested positive, nursing home
06	new case, tested positive, wear mask, president trump	new case, confirmed case, report new, single day, juen juen
07	tested positive, wear mask, president trump, new case	new case, pandemic began, case today, new death

**Table 5.3.1** Trending bigrams on rumours and non-rumours over months flow

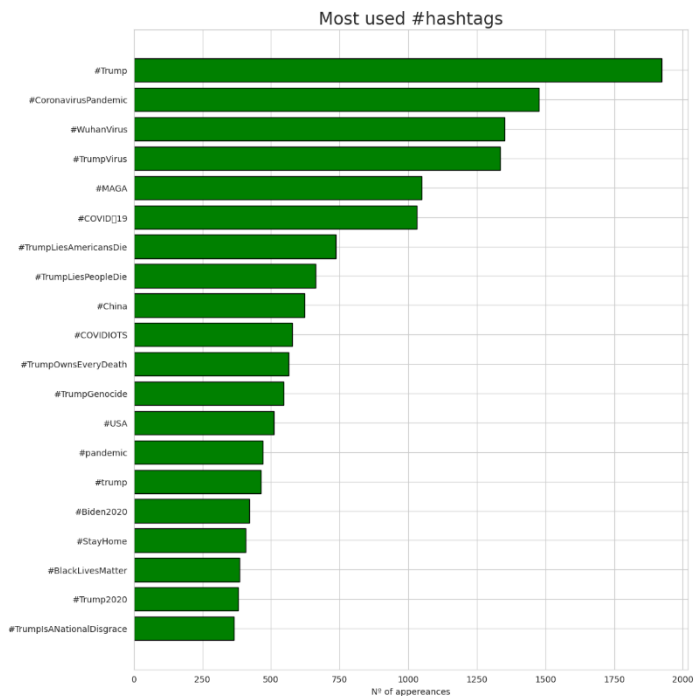


Figure 5.3.9. non-rumours top 20 popular hashtags.

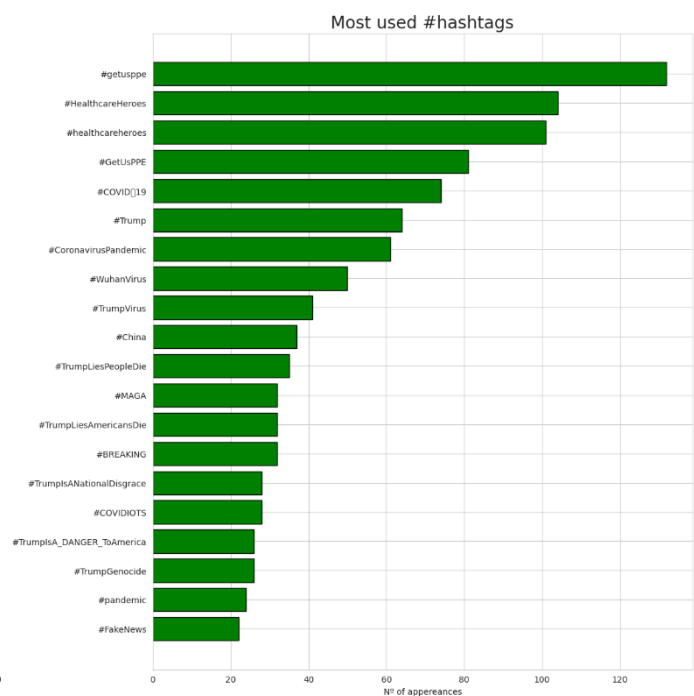


Figure 5.3.10. rumours top 20 popular hashtags.

	non-rumours	rumours
#Follower	5409719	4776971
#Following	7931	7196
Ratio (#Follower / #Following)	682	663
#Favrate Counts	14299	5500
#Post	112456	151732
Account Created Date	2011-04-30	2010-12-01
Geo Enabled	42%	40%

**Table 5.3.2** User Statistics

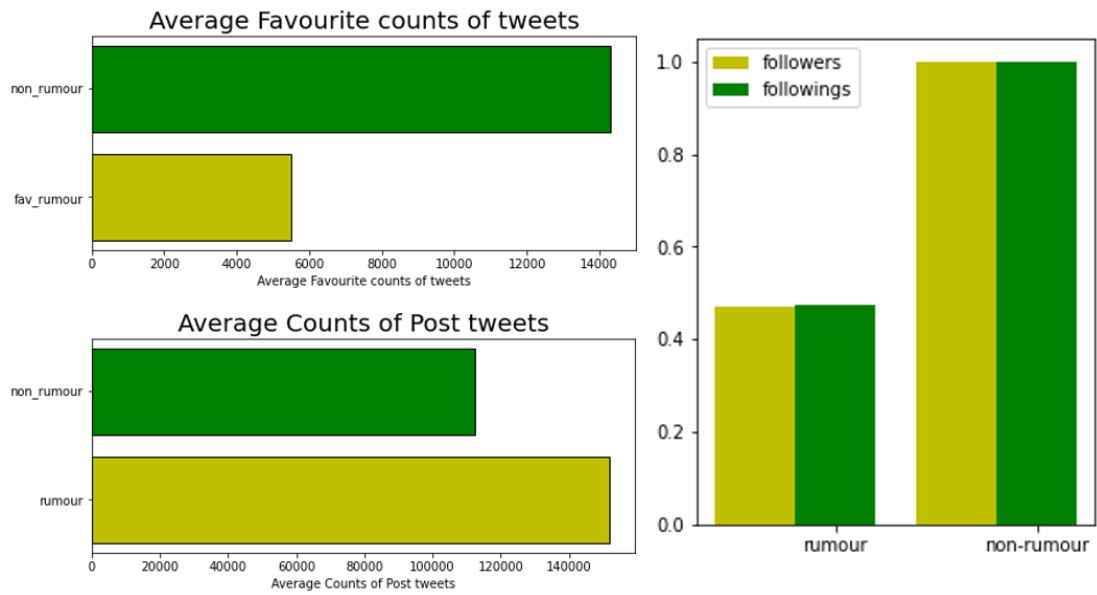


Figure 5.3.11. user statistic visualization

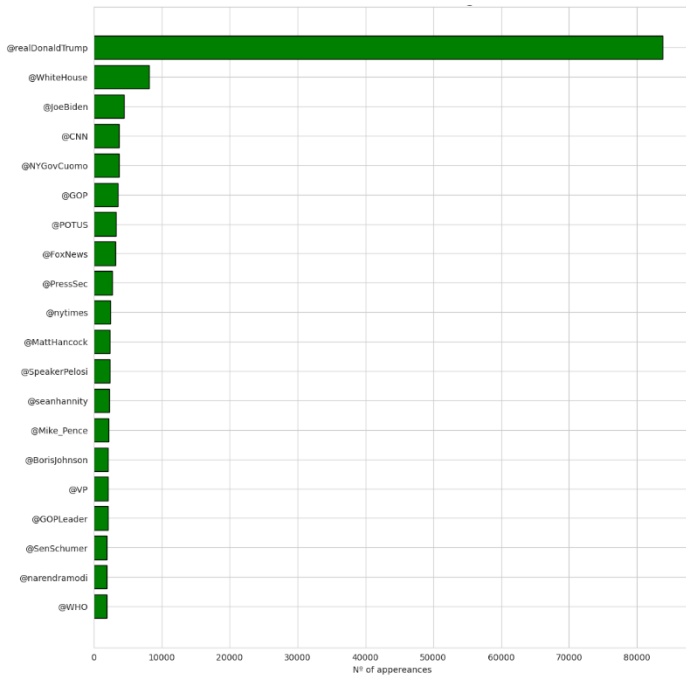


Figure 5.3.12. most active @users in non-rumours.

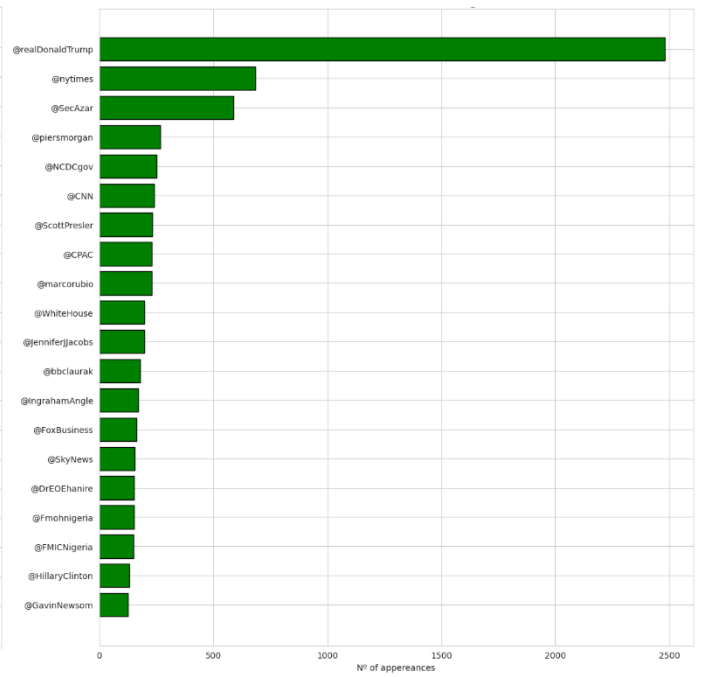
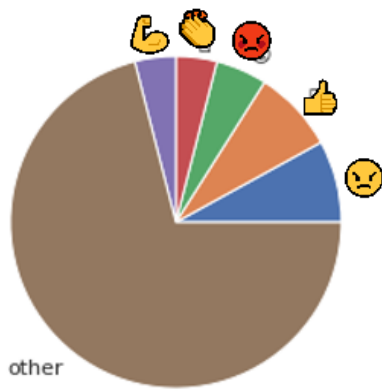
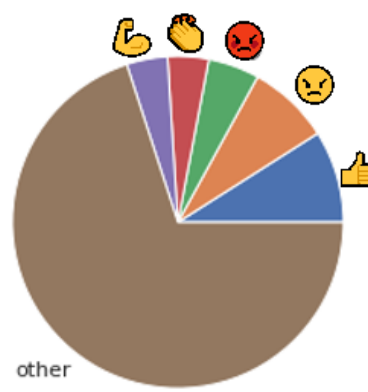


Figure 5.3.13. most active @user in rumours.



**Non-rumour emotions**

Figure 5.3.13. (a) non-rumour emoji emotions



**rumour emotions**

(b) rumour emoji emotions

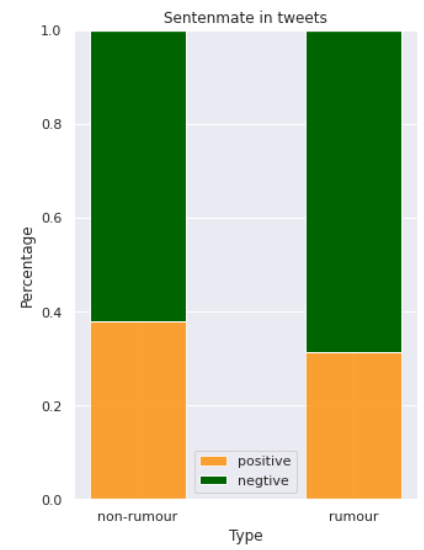


Figure 5.3.14. Sentiments